## CLAIMS

1     1. A method of detecting duplicates in a set of documents having associated
2     nearest neighbor similarity scores, the method including:

3         for a particular document in the set of documents, selecting nearest neighbors of
4         the particular document; and

5         flagging as potential duplicates the nearest neighbors of the particular document
6         that have respective nearest neighbor similarity scores that are identical.

1     2. The method of claim 1, further including flagging as potential duplicates the
2     nearest neighbors of the particular document that have respective nearest neighbor
3     similarity scores that are within a tolerance t of one another.

1     3. The method of claim 1, wherein the nearest neighbor similarity scores are
2     calculated prior to duplicate detection for a different purpose than the duplicate
3     detection and stored with the documents.

1     4. The method of claim 1, wherein the k nearest neighbors are determined prior
2     to duplicate detection for a different purpose than the duplicate detection and stored
3     with the documents.

1     5. The method of claim 1, wherein the documents are text documents.

1     6. The method of claim 5, wherein the text documents include visual formatting.

1     7. The method of claim 1, wherein the documents are voice recordings.

1     8. The method of claim 1, wherein the documents are musical performances.

1     9. The method of claim 1, wherein the documents are graphic images.

1     10. The method of claim 1, wherein the nearest neighbors are limited to k nearest
2     neighbors.

1     11. A method of detecting duplicates in a set of documents, the method including:

2     identifying nearest neighbors of documents in the set of documents, based on

3     nearest neighbor similarity scores;

4     for a particular document in the set of documents, flagging as potential duplicates

5     the nearest neighbors of the particular document that have respective nearest

6     neighbor similarity scores that are identical.

1     12. The method of claim 11, further including flagging as potential duplicates the

2     nearest neighbors of the particular document that have respective nearest neighbor

3     similarity scores that are within a tolerance t of one another.